

Corpus Building with TEC Tools Tutorial

Sally Marshall, Sofia Malmatidou and Kyung-Hye Kim
University of Manchester

| Contents | Page |
|------------------------------------------------------|-------------|
| 1. Scanning and Converting Images to Text | 2 |
| 1.1 Scanning Documents | |
| 1.2 Choosing OCR Software | |
| 1.3 Extracting Text | |
| 2. Preparing Text and Header Files | 8 |
| 2.1 Filenames | |
| 2.2 Text Files and Header Files | |
| 2.3 DTD Files | |
| 2.4 Setting up jEdit | |
| 2.5 Preparing Text Files | |
| 2.6 Preparing Header Files | |
| 3. Building a Corpus | 16 |
| 3.1 Introduction to TEC Tools | |
| 3.2 Installing the Indexer | |
| 3.3 Adding Data | |
| 3.4 Setting up the Indexer | |
| 3.5 Indexing Files | |
| 3.6 Testing the Corpus | |
| 3.7 Sharing a Corpus | |
| 4. Searching a Corpus and Saving Concordances | 24 |
| 4.1 The Corpus Browser Interface | |
| 4.2 Searching a Corpus | |
| 4.3 Sorting Concordances | |
| 4.4 Saving Search Results | |
| 5. Other Resources | 30 |

Contributors:

Kyung-Hye Kim (University of Manchester)

Sofia Malmatidou (University of Manchester)

Sally Marshall (University of Manchester)

Editor:

Sally Marshall (University of Manchester)

1. Scanning and Converting Images to Text

Once you have selected the texts you plan to include in your corpus, and obtained permission from the copyright holder (if necessary), you will need to convert your texts into a computer-readable, or digital/electronic, format.

If your data is on printed paper, you will need to scan the printed documents and use OCR (Optical Character Recognition) Software to extract text from them. If you are using texts that are already computer-readable (e.g. text taken from the internet, e-Books, searchable PDF files, Word files etc) you can move on to step2.

1.1 Scanning Documents

Depending on the type and volume of printed material you are working with you may choose to use either a flat-bed scanner or a scanner with a document feeder.

Flat-bed scanner

Place one page at a time onto the glass scanner bed

Most home-office scanners are flat-bed scanners

Scanner with document feeder

Place multiple pages in the document feeder

Many print shops and university departments have scanners with document feeders

If you have a large number of books to scan you may want to consider removing the pages from their binding so you can use a document feeder. Many print shops are able to remove pages from book bindings using a guillotine, and may also be able to rebind books after scanning. When you scan multi-page documents you should save all pages as one PDF file.

Whichever type of scanner you decide to use you should bear in mind that the quality of the scanned images should be as high as possible in order to ensure good results when using OCR software. Therefore you should check that your scanner is capable of producing images of a sufficiently high quality. Some other points to consider include:

Paper quality

Thicker paper is less prone to “show-through” (where text from the other side of the page shows through). For this reason, hard-backed editions of books may produce better results than paperbacks.

Print quality

The better the quality of the printed text, the better the quality of the output. Avoid poor quality print, for example with “bleed”.

Font

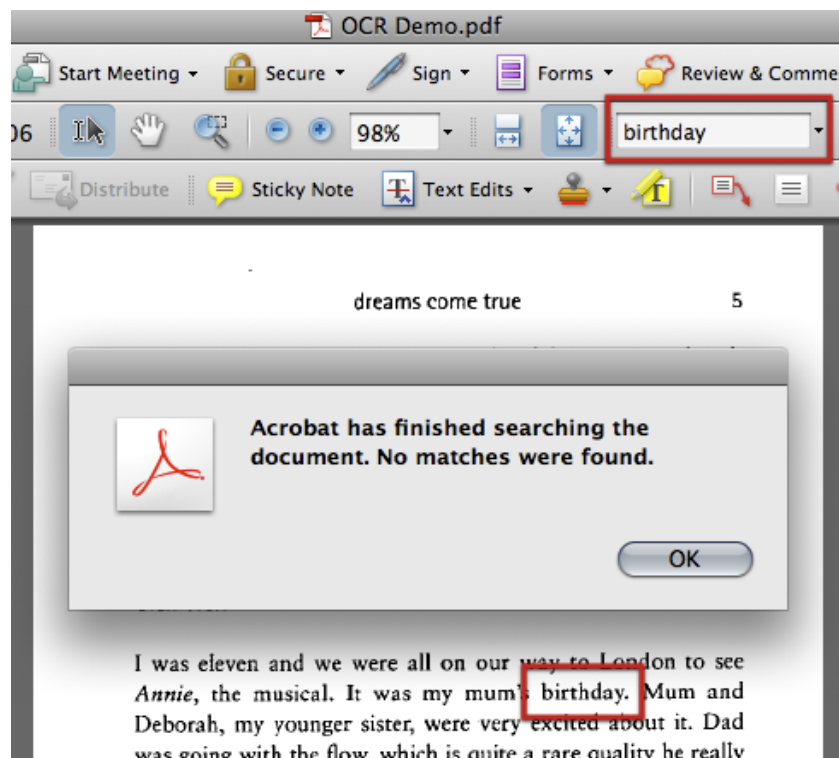
Commonly used fonts often produce better results than unusual fonts

Scanner settings

You should set your scanner to a minimum setting of 300dpi, and check for other settings that may affect the quality of images.

After scanning your printed materials you will have a number of PDF image files. These files are - like photographs - copied images of the printed material and cannot be searched for text. If you open one of these PDF files and try to use the Find function, it will not return any results.

For example,



In order to convert PDF image files to searchable text, you will need to perform OCR.

1.2 Choosing OCR Software

In order to carry out OCR on your PDF images you will need to choose and install OCR software. Many software packages have demo or trial versions available for download, and you might also want to check whether your scanner came with any OCR software included.

Depending on the language you wish to scan, and the operating system you use (Windows or Mac), there are a number of different OCR software packages available. Some of these are listed at the end of this section, but you are advised to search online for the most up-to-date OCR software, as OCR technology (especially for non-

alphabet languages) is continually improving. With good-quality PDF images, English OCR software is able to recognize text with a very high degree of accuracy.

Particularly if you are using non-alphabet languages, you are advised to test a few OCR software packages and compare their performance. If you carry out the OCR process on a sample document using a number of different software packages, you can compare accuracy by searching for a number of words to see which software has been most successful in recognizing them. You can also export the text to see which software preserves the original format best.

Example

Four OCR software packages were tested with Japanese data. The OCR output accuracy was compared, then the exported text formatting of the top-performing two software packages was compared.

| OCR Software | Test 1 Accuracy Ranking | Test 2 Formatting Ranking |
|--------------|-------------------------|---------------------------|
| A | 1= | 2 |
| B | 3 | |
| C | 2 | |
| D | 1= | 1 |

Software A and D performed equally well for character recognition accuracy, but software D was able to maintain the original text formatting more successfully. On the basis of the two phases of comparison, software D was selected for use.

*****UNDER CONSTRUCTION*****

| OCR Software | Operating System | | Language | | | | | | |
|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| | Windows | Mac | English | European | Greek, Russian | Chinese | Japanese | Korean | Arabic |
| Abbyy Finereader | <input type="checkbox"/> | | | | | | <input type="checkbox"/> | | |
| AbbyyPDF Transformer | | | | | | | | | |
| Adobe Acrobat Pro | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| NuanceOmniPage Pro | | | | | | | | | |
| IrisReadIris | | | | | | | | | |
| | | | | | | | | | |
| PanasonicYomitoriKakumei | <input type="checkbox"/> | | <input type="checkbox"/> | | | | <input type="checkbox"/> | | |
| EpsomYonde Koko | <input type="checkbox"/> | | <input type="checkbox"/> | | | | <input type="checkbox"/> | | |
| KodenshaArumi | <input type="checkbox"/> | | | | | | | <input type="checkbox"/> | |
| | | | | | | | | | |
| | | | | | | | | | |
| | | | | | | | | | |

*****UNDER CONSTRUCTION*****

NB:

This list is not exhaustive and is intended only as a sample. You are advised to carry out your own research to find suitable OCR software.

You should check the system requirements for each software package for compatibility. For example, some OCR software can only be used with localized editions of Windows.

Mac users currently have fewer OCR software options available and so may want to consider using a Windows PC to do the OCR process. Alternatively you could use Parallels or Boot Camp on your(Intel) Mac to install Windows.

1.3 Extracting Text

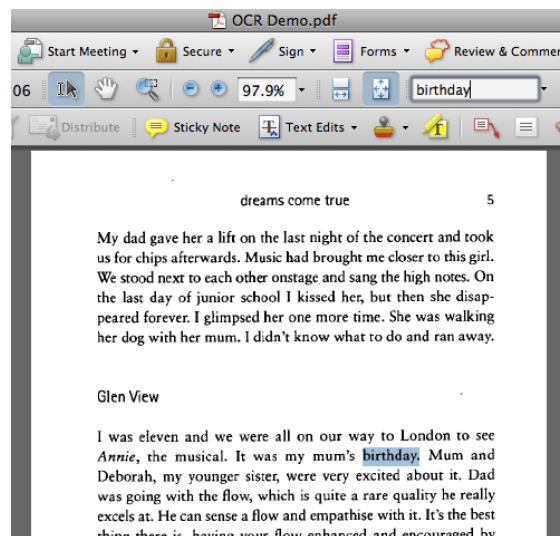
Once you have selected the OCR software you will use, you can carry out OCR on your PDF images.

Different OCR software packages have different settings and options. You should familiarize yourself with various settings to achieve the best results. If there is an option for “Autocorrect for skew” (which corrects minor rotations to pages after scanning) you should check this. Make sure you have set the language correctly before carrying out OCR.

When doing OCR on documents with complicated layouts (such as newspapers or scientific articles), consider using the Crop tool to select sections, and perform OCR section by section. This can help avoid formatting problems.

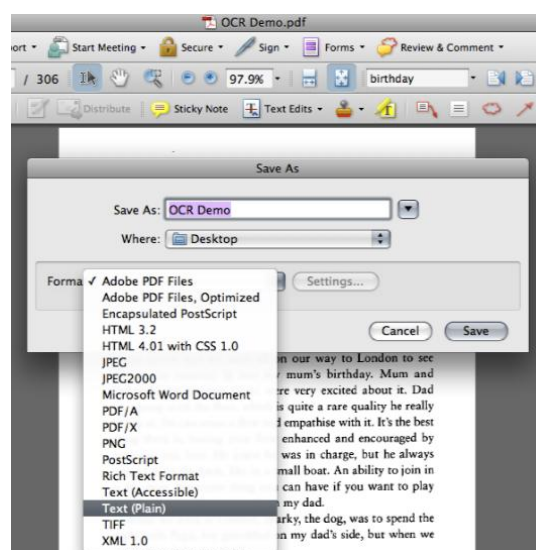
Once you have carried out OCR on your PDF images you will be able to search them for text using the Find function. Try searching for a word that you know appears in the text. When you save a copy of the PDF file after carrying out OCR, the PDF is now **searchable**.

e.g.

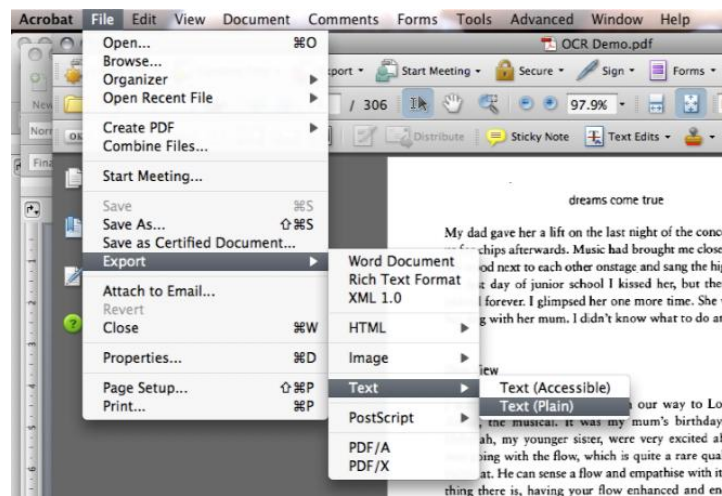


In order to carry out the next stage of the corpus-building process you need to use **searchable and editable** text. To do this, you should look for options in your OCR software that allow you to either:

(1) **Save As** a Word or text document



or, (2) **Export** the data as text.



After you have saved the text as a separate file (.doc or .txt or similar) you will probably have three files for each of your texts:

- a PDF image
- a searchable, OCRd PDF
- a searchable and editable text file

For the next stage of the process, you will need to use these text files, which are now your raw data.

2. Preparing Text and Header Files

This section will help you turn your raw data into text files that are suitable for use with the TEC Corpus Software.

2.1 Filenames

You will need to consider what system of filenames is most useful for your purposes. The conventions you adopt for naming files in your corpus will probably be related to the design of your corpus.

It is possible to identify types of texts (i.e. sub-corpora) using a system of filenames that indicates the text type. For example, in the TEC corpus fictional texts all share the letters **fn** in their filename, and biographical texts share **bb**.

Since the TEC Tools software was originally designed for use with a corpus comprising only texts translated into English, you may wish to devise a system of filenames that helps you identify the source and target language of your translated texts, and/or differentiate between translated and non-translated texts.

Example

If a corpus includes English texts and their Japanese translations, and Japanese texts and their English translations there will be four types of texts. In such a case, texts can be given filenames such as:

| | | | |
|-------------------|------|----------------------|------|
| English original | EE01 | Japanese translation | EJ01 |
| Japanese original | JJ02 | English translation | JE02 |

These filenames allow the identification of sub-corpora - translated vs. non-translated texts (determined by whether the letters denoting language are the same or not), Japanese or English texts (the second letter is the language of the text) - and identify ST-TT pairs (the filename number identifies pairs of texts).

e.g.

| | |
|----------------------------|-----------------------------|
| Filename EJ07 | Filename JJ04 |
| Translated text | Non-translated text |
| Source Language = English | Language = Japanese |
| Target Language = Japanese | |
| Source Text Filename= EE07 | Target Text Filename = JE04 |

Other systems of filenames can also be developed to encode various relevant attributes.

2.2 Text Files and Header Files

A corpus is typically made up of text files and header files.

- **Text files** contain the actual data to be analysed.
- **Header files** contain meta-data such as the title of the text, author, publisher and any other features of interest.

Text files and header files come in pairs, and the filenames for text files and their associated header files should be the same. Text files can be identified by a **.xml** extension and header files have a **.hed** extension.

e.g.

| Text File | Header File |
|----------------|----------------|
| bb00000001.xml | bb00000001.hed |

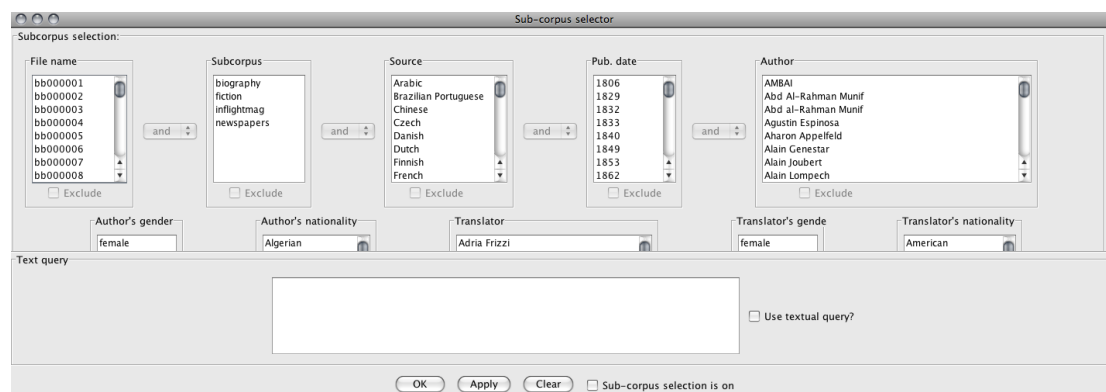
The TEC corpus and many other corpora use such a system of text files and associated header files because it allows the user to select sub-corpora - i.e. sets of texts that share certain features - to search in the corpus Browser.

Below is an example of a header file.

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<!DOCTYPE techheader SYSTEM "techheader.dtd">
<techheader>
  <title subcorpusid="biography" filename="bb000010">
    <subcorpus>biography</subcorpus>
    <collection>Girls of Alexandria</collection>
  </title>
  <section id="s1">
    <translator gender="female">
      <name>Frances Liardet</name>
      <nationality description="British"/>
      <employment>writer, translator</employment>
      <status>freelance, part-time</status>
    </translator>
    <translation extent="61739">
      <publisher>Quartet Books</publisher>
      <pubPlace>UK</pubPlace>
      <date year="1993">1993</date>
      <copyright>Frances Liardet</copyright>
    </translation>
    <translationProcess mode="wst">
      <direction>into mother tongue</direction>
      <type>full</type>
    </translationProcess>
    <author gender="male">
      <name>Edwar al-Kharrat</name>
      <nationality description="Egyptian"/>
    </author>
    <sourceText>
      <language>Arabic</language>
      <status>original</status>
    </sourceText>
  </section>
</techheader>
```

You can see from this example that various pieces of information relevant to the text are listed. These include the name, nationality and gender of the author and translator, and the title, publisher and language of the text.

This information can be used to define sub-corpora when browsing the TEC corpus, as you can see from this screen-shot of the Sub-corpus Selector window of the Corpus Browser interface.



Some of the information in the header file can be input freely, for example, the name of the author or the title of the book. However, some information must be selected from a pre-determined list of options, for example, gender must be either male or female.

The rules that determine what can and cannot be included in a header file are defined in a file called a **DTD file**.

2.3 DTD Files

The rules that define what constitutes a “well-formed” text file or header file are listed in DTD files. When the TEC Tools software processes your text files and header files it needs to have access to an accompanying DTD file for the text files and another DTD file for the header files.

When you download the TEC Tools software, two DTD files are included in the folder: tectext.dtd and techeader.dtd.

After you have prepared all of your text files you will need to save them in a folder that also contains the text DTD file. Similarly, all of your header files must be saved in a separate folder that also contains the header DTD file.

2.4 Setting up jEdit

In order to prepare both text and header files, you will need to use text editing software that can be used for XML encoding. jEdit is a freely available text editor that works on both Windows and Mac operating systems, among others. To start

using jEdit, download and install the software (<http://www.jedit.org/>). The first time you use jEdit you will need to set it up. To do this, follow the following steps:

- Open jEdit
- Go to Utilities > Global Options > Encoding
- Select "Default Character Encoding UTF-8"
- Deselect "Auto-detect file encoding where possible"
- Go to Utilities > Global Options > Editing>Word Wrap
- Select "Soft".
- Click on Apply and OK.
- Go to Plugins > Plugin Manager> Install >XML
- Click Install
- Close
- Go to Plugins > Sidekick
- Check "Parse on keystroke"
- Check "Highlight markers in structure browsers"

By doing this you are preparing jEdit to work with XML encoding, and enabling jEdit to check for Errors in your files.

2.5 Preparing Text Files

Although there is considerable variation in the time it takes to prepare text files, based on the quality of the OCR text output and other variables, you might expect to spend between 1-2 hours to prepare the text from a 200-page book.

- Open your text file in Word.
- Do a spellcheck.
- Notice any repeated OCR errors (e.g. "Tm" for "I'm" etc) and correct these using Find and Replace.

As the text files will later be encoded with XML encoding - which enables the corpus software to process the text - you will also need to remove certain characters from your text, as they have specific meanings in XML.

Using Find and Replace,

- Find all instances of "<" and Replace them with nothing.
- Find all instances of ">" and Replace them with nothing.
- Find all instances of "&" and Replace them with "and".

Then, **save the file as a text file with UTF-8 encoding.**

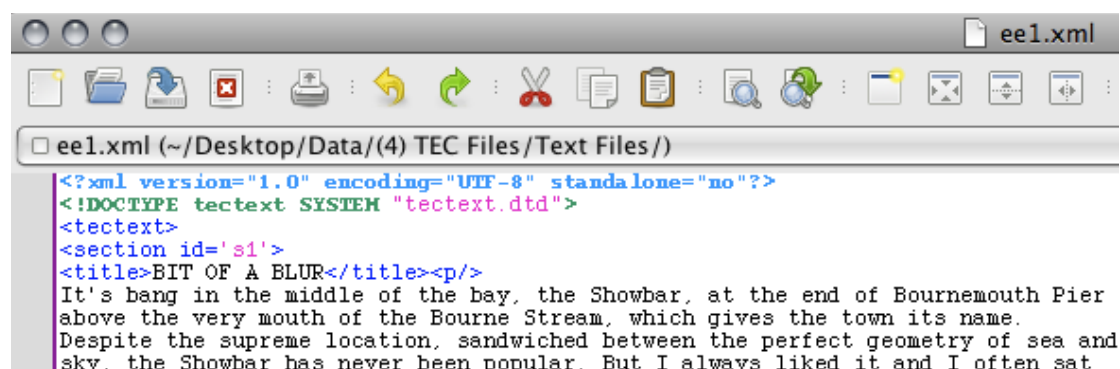
There are many ways to tag your text with information, and you should consider what information you will need in order to analyse your corpus. The most basic kind of XML tagging for text files is described below.

- Open the UTF-8 text file in jEdit.
- Remove any white spaces (i.e. delete blank lines)
- Either delete parts of the text that you do not wish to analyse (e.g. chapter headings, footnotes etc) or use <omit> tags to tell the Corpus Browser not to include this material in searches.
- Add the following XML tagging at the start of the body of text:

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<!DOCTYPE tectext SYSTEM "tectext.dtd">
<tectext>
<section id='s1'>
<title>INSERT TITLE</title><p/>
```

- Replace INSERT TITLE with the title of your text.

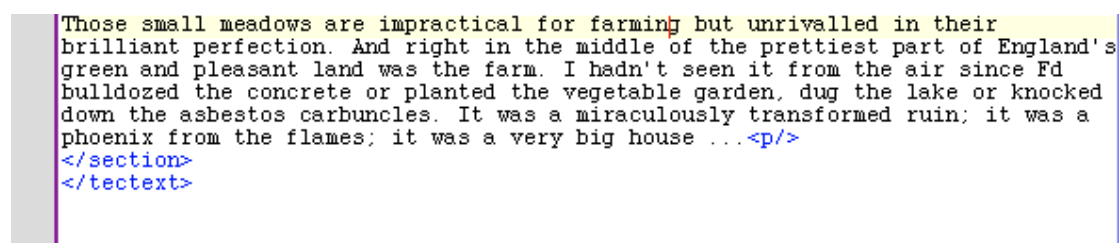
e.g.



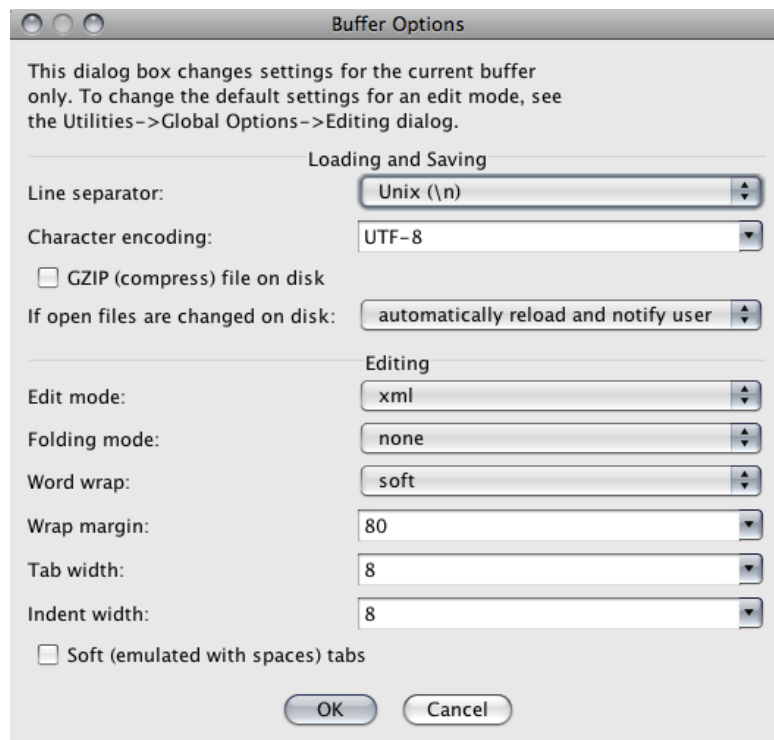
- Add the following XML tagging at the end of the text:

```
<p/>
</section>
</tectext>
```

e.g.



- Go to Utilities > Buffer Options. Use the following settings:



- Save and close jEdit.
- Put your text file in a folder with the text DTD file.
- Reopen the text file in jEdit.
- Check for errors by going to Plugins > Error List >Errorlist.

There are several ways to prepare texts and you may find that you use an alternative method to the one described here. Whichever method you use, you should be consistent in your approach and may find it useful to create a checklist like the one below.

NB. The checklist is an **example** only.

Sample Checklist for Preparing Text Files

| | | | | | | | | | | | | | | | | | |
|--------------------------------------------------------------|---------------------------------------|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
| Filename | e.g. ej01 | | | | | | | | | | | | | | | | |
| (1) Adobe Acrobat | | | | | | | | | | | | | | | | | |
| Export text | | | | | | | | | | | | | | | | | |
| (2) MS Word | | | | | | | | | | | | | | | | | |
| Spellcheck text | | | | | | | | | | | | | | | | | |
| Save as UTF-8 | | | | | | | | | | | | | | | | | |
| (3) jEdit – Text Preparation | | | | | | | | | | | | | | | | | |
| XML | Find & replace < with nothing | | | | | | | | | | | | | | | | |
| | Find & replace > with nothing | | | | | | | | | | | | | | | | |
| | Find & replace & with “and” | | | | | | | | | | | | | | | | |
| Check for Repeating OCR Errors | Find & replace ^ with nothing | | | | | | | | | | | | | | | | |
| | Find & replace _1_ with “_I_” | | | | | | | | | | | | | | | | |
| | Find & replace _!_ with “_I_” | | | | | | | | | | | | | | | | |
| | Find & replace “TTi” with “Th” | | | | | | | | | | | | | | | | |
| | Find & replace “TU_” with “I’ll_” | | | | | | | | | | | | | | | | |
| | Find & replace “Td_” with “I’d_” | | | | | | | | | | | | | | | | |
| | Find & replace “Tm_” with “I’m_” | | | | | | | | | | | | | | | | |
| | Find & replace “T_” with “I_” | | | | | | | | | | | | | | | | |
| | Find & replace “iH” with “ill” | | | | | | | | | | | | | | | | |
| | etc ... | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |
| Simultaneously | Omit front matter up to start of text | | | | | | | | | | | | | | | | |
| | Omit chapter headings, titles etc | | | | | | | | | | | | | | | | |
| | Remove blank lines | | | | | | | | | | | | | | | | |
| | Check for other errors | | | | | | | | | | | | | | | | |
| | Repair broken lines | | | | | | | | | | | | | | | | |
| | etc ... | | | | | | | | | | | | | | | | |
| | Save regularly (as bookname.txt) | | | | | | | | | | | | | | | | |
| (4) jEdit – XML Tagging | | | | | | | | | | | | | | | | | |
| Insert title tagging from sample file | | | | | | | | | | | | | | | | | |
| Change title name | | | | | | | | | | | | | | | | | |
| Insert end tagging from sample file | | | | | | | | | | | | | | | | | |
| Set Utilities> Buffer Options > Line Separator > “Unix (\n)” | | | | | | | | | | | | | | | | | |
| Save As (filename.xml) | | | | | | | | | | | | | | | | | |
| Close jEdit | | | | | | | | | | | | | | | | | |
| Put in folder with tectext.dtd | | | | | | | | | | | | | | | | | |
| Open the file | | | | | | | | | | | | | | | | | |
| Check for errors: Plugins > Error List | | | | | | | | | | | | | | | | | |

2.6 Preparing Header Files

Currently being written...

3. Corpus Building with TEC Tools

This part of the tutorial introduces the process of downloading and using the TEC Tools corpus software (also known as the MODNLP corpus suite).

3.1 Introduction to TEC Tools

The TEC Tools corpus software is a set of corpus tools that was developed by Dr Saturnino Luz for use with the Translational English Corpus (<http://ronaldo.cs.tcd.ie/tec2/jnlp/>) and designed to allow free access to linguistic material over the internet (Luz, 2011). However, the software can also be downloaded and used to create and search a corpus using any texts you like. Originally the tools supported only English (and other European languages) but are currently being developed to allow use with non-alphabet languages such as Japanese¹.

The software consists of three modules:

- An **indexer** (called modnlp-idx) which allows you to create an index.
- A **corpus browser** (called modnlp-teccli) which can be used to select sub-corpora (by accessing indexes) and browse concordances.
- A **corpus server** (called modnlp-tecser) which can be used to make data and concordances (although not necessarily full texts) available to other users over the internet.

System Requirements

- The MODNLP software works with both Windows and Mac operating systems.
- You should ensure that your computer is using the most up-to-date version of Java available for its operating system.

3.2 Installing the Indexer

- Download the Indexer (IDX) which is available here:
<http://ronaldo.cs.tcd.ie/~luzs/tmp/modnlp-idx-0.2.0-bin-jp.tgz>
- Unzip the folder to access the following files:

¹ These instructions apply to the developer 'multilingual' version currently being tested (Jan 2011). Please check for the latest version of the software before starting.

| Name | Date Modified | Size | Kind |
|------------------------|-----------------------|--------|---------------|
| antlr-2.7.6.jar | 4 January 2011, 15:51 | 404 KB | Java JAR file |
| commons-pool-1.2.jar | 4 January 2011, 15:51 | 48 KB | Java JAR file |
| COPYING-libs | 4 January 2011, 15:51 | 68 KB | Plain text |
| ep | Today, 15:05 | -- | Folder |
| eph | Today, 15:05 | -- | Folder |
| exist-modules.jar | 4 January 2011, 15:51 | 72 KB | Java JAR file |
| exist.jar | 4 January 2011, 15:51 | 3.3 MB | Java JAR file |
| gnu-regexp.jar | 4 January 2011, 15:51 | 36 KB | Java JAR file |
| idx.jar | 4 January 2011, 15:51 | 144 KB | Java JAR file |
| idxmgr.properties | Today, 12:33 | 4 KB | TextE...ment |
| je.jar | 4 January 2011, 15:51 | 2 MB | Java JAR file |
| jgroups-all.jar | 4 January 2011, 15:51 | 1.6 MB | Java JAR file |
| jp | 1 January 2011, 00:11 | -- | Folder |
| jph | 4 January 2011, 15:14 | -- | Folder |
| log4j-1.2.14.jar | 4 January 2011, 15:51 | 392 KB | Java JAR file |
| README | 4 January 2011, 15:51 | 4 KB | Plain text |
| resolver.jar | 4 January 2011, 15:51 | 68 KB | Java JAR file |
| runidx.sh | 4 January 2011, 15:51 | 4 KB | Plain text |
| sunxacml.jar | 4 January 2011, 15:51 | 264 KB | Java JAR file |
| TUTORIAL.txt | 4 January 2011, 15:51 | 8 KB | Plain text |
| xmldb.jar | 4 January 2011, 15:51 | 16 KB | Java JAR file |
| xmlrpc-1.2-patched.jar | 4 January 2011, 15:51 | 116 KB | Java JAR file |

The **TUTORIAL.txt** file, highlighted purple, is the document on which the current instructions are based.

Other files and folders that you need to use are highlighted in grey.

- The **idx.jar** file is the **indexer**.
- The **idxmgr.properties** is a file that contains **settings** that are used by the indexer.
- The folders **ep** and **eph** are for storing English **data** files.
- The folders **jp** and **jph** are for storing Japanese **data** files.

3.3 Adding Data

- In the **modnlp-idx-0.2.0-bin-tec** folder, create a folder called **data**.

(a) Adding English-language data

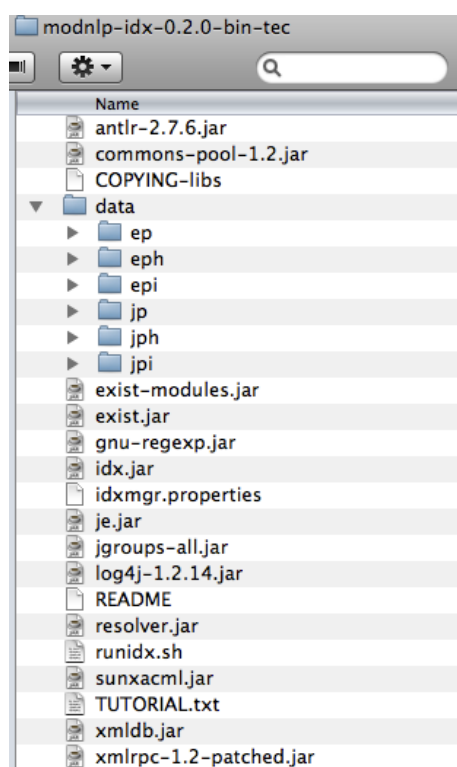
- Move the folders **ep** and **eph** into this **data** folder.
- When you have prepared your **text files** (.xml files), put them into the **ep** folder, which already contains a dtd file called **ep.dtd** (and a couple of sample text files).
- When you have prepared your **header files** (.hed files), put them into the **eph** folder which already contains a dtd file called **ep.dtd** (and a couple of sample header files).
- Create another folder called **epi** and put this into the **data** folder too. This will be used later to store index files.

(b) Adding Japanese-language data

- Move the folders **jp** and **jph** into the **data** folder.

- When you have prepared your **text files** (.xml files), put them in the **jp** folder which already contains a dtd file called **tectext.dtd** (and a couple of sample text files).
- When you have prepared your **header files** (.hed files), put them in the **jph** folder which already contains a dtd file called **techeader.dtd** (and a couple of sample header files).
- Create another folder called **jpi** and put this into the **data** folder. This will be used later to store index files.

If you have added both English and Japanese data, you should have the following sub-folders in your data folder:

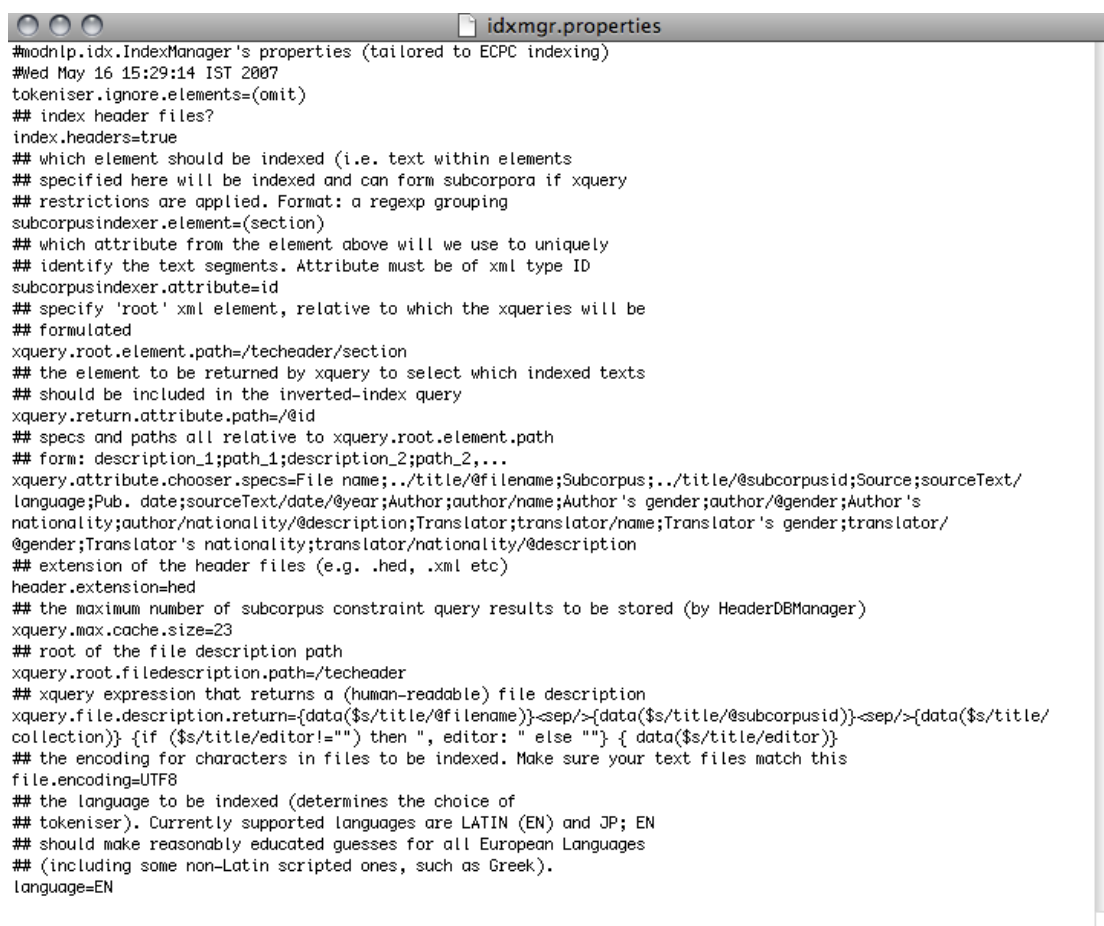


NB. The data folders **ep** and **eph** (**jp** and **jph**) already contain **.dtd** files and some sample text and header files. You may prefer not to replace these with your own data until you have completed the installation and testing process. If you have successfully tested the Corpus Browser using the sample data provided, you will be familiar with the procedure. If you then encounter problems when using your own data, troubleshooting may be more straightforward. When you have added your own data files to the data folders, remember to remove the **sample** files that were there from the start, but leave the **.dtd** files.

3.4 Setting up in the Indexer

idxmgr.properties tells the index manager how to tokenise text files, which sections of the text to index, and which parts of the corpus should be included in a sub-corpus etc.

The idxmgr.properties file can be opened using any text editor.



```
#modnlp.idx.IndexManager's properties (tailored to ECPC indexing)
#Wed May 16 15:29:14 IST 2007
tokeniser.ignore.elements=(omit)
## index header files?
index.headers=true
## which element should be indexed (i.e. text within elements
## specified here will be indexed and can form subcorpora if xquery
## restrictions are applied. Format: a regexp grouping
subcorpusindexer.element=(section)
## which attribute from the element above will we use to uniquely
## identify the text segments. Attribute must be of xml type ID
subcorpusindexer.attribute=id
## specify 'root' xml element, relative to which the xqueries will be
## formulated
xquery.root.element.path=/techeader/section
## the element to be returned by xquery to select which indexed texts
## should be included in the inverted-index query
xquery.return.attribute.path=/@id
## specs and paths all relative to xquery.root.element.path
## form: description_1;path_1;description_2;path_2,...
xquery.attribute.chooser.specs=File name;../title/@filename;Subcorpus;../title/@subcorpusid;Source;sourceText/
language;Pub. date;sourceText/date/@year;Author;author/name;Author's gender;author/@gender;Author's
nationality;author/nationality/@description;Translator;translator/name;Translator's gender;translator/
@gender;Translator's nationality;translator/nationality/@description
## extension of the header files (e.g. .hed, .xml etc)
header.extension=hed
## the maximum number of subcorpus constraint query results to be stored (by HeaderDBManager)
xquery.max.cache.size=23
## root of the file description path
xquery.root.filedescription.path=/techeader
## xquery expression that returns a (human-readable) file description
xquery.file.description.return={data($s/title/@filename)}<sep/>{data($s/title/@subcorpusid)}<sep/>{data($s/title/
collection)} {if ($s/title/editor!="") then ", editor: " else ""} { data($s/title/editor)}
## the encoding for characters in files to be indexed. Make sure your text files match this
file.encoding=UTF8
## the language to be indexed (determines the choice of
## tokeniser). Currently supported languages are LATIN (EN) and JP; EN
## should make reasonably educated guesses for all European Languages
## (including some non-Latin scripted ones, such as Greek).
language=EN
```

In this version of the software it is possible to set the indexer to work with English (EN) or Japanese (JP) text files that are encoded in UTF-8.

When the setting is EN the indexer should also be able to handle European languages. (See screenshot bottom 5 lines.)

If you want to build a corpus that contains both English texts and Japanese texts, you will need to change the **language setting** between JP and EN when indexing each set of data files. You should complete the indexing process for your English-language data, and then repeat for your Japanese-language data. When you carry out the indexing process, you should check the language setting is correct for your current data.

To change the language setting:

- OPEN idxmgr.properties (using TextEdit or any other text editor)
- Go to the bottom line that says 'language=JP' or 'language=EN'.
- Change the setting to 'language=EN' or 'language=JP' as you require.
- SAVE and CLOSE

3.5 Indexing Files

- Open **idx.jar**. (NB. It may appear as just **idx**).
- A window will appear asking you to “Choose a location for the index”. This location is where the Indexer will store the index it creates. The index will later be used by the Corpus Browser.
- In this window, open the **data** folder and select (i.e. highlight but do not open) either sub-folder **epi** or **jpi** (depending what set of data you are currently indexing).
- Click “Choose a location for the index”.
- Next you will be asked to choose the folder where header files are stored. Select **epi** or **jpi**.
- After about 30 seconds pause you will be asked to give a URL for public access to the header files. This is relevant if you want to make your corpus available to other users. Click “OK” to accept the default URL.
- The main indexer window should now appear. Click “Index New Files”. Open the **epi** or **jpi** folder and select the XML files you want to index. Click “Index New Files”.
- When the indexing process is finished, click QUIT > Yes.

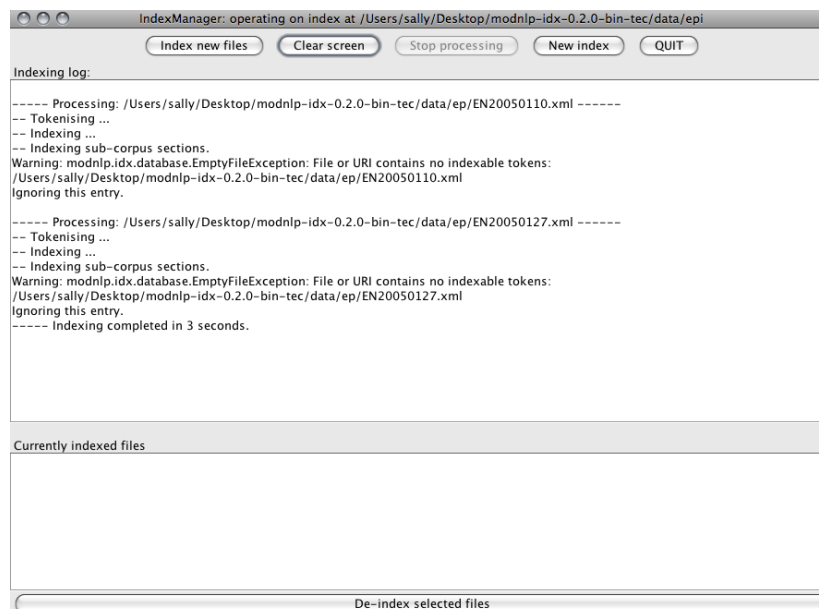
The indexer will have stored an index in the **epi** or **jpi** folder.

Troubleshooting:

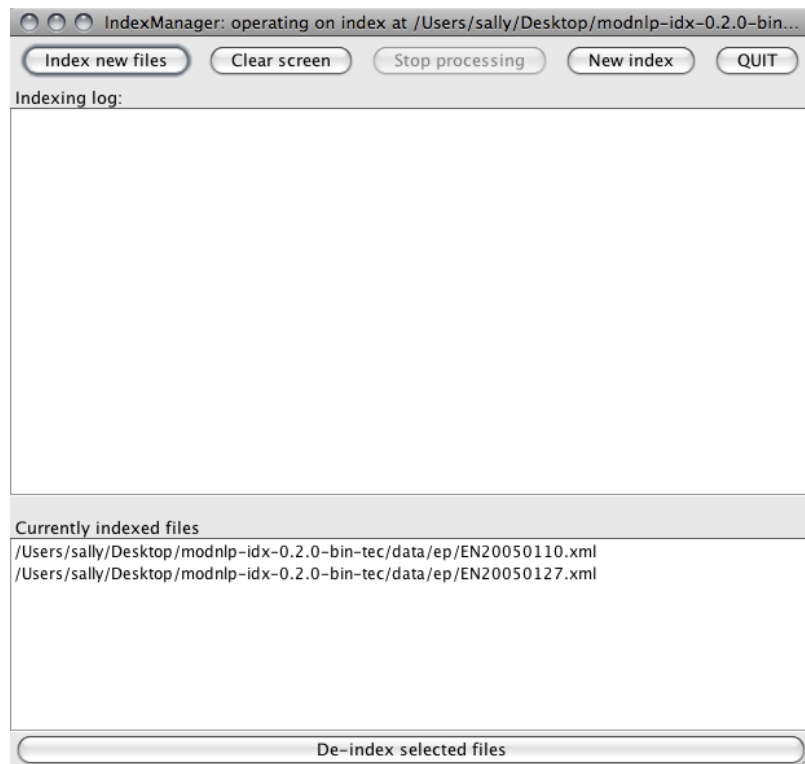
Check that you have the correct language setting in the **idxmgr.properties**.

Check that you *selected* but did not *open* folders.

If you get an error message such as this:



Close the Indexer and reopen it. You may find that the issue has been resolved.



3.6 Testing the Corpus

To search your corpus you will need to use the Corpus Browser (TECCLI). The Corpus Browser works with the index created by the Indexer (IDX) as described above to select and search your texts.

- Download the Corpus Browser (TECCLI) which is available here:
<http://ronaldo.cs.tcd.ie/~luzs/tmp/modnlp-teccli-0.7.0-bin-tec.tgz>
- Unzip the folder to access the following files:

| modnlp-teccli-0.7.0-bin-tec | | | | |
|-----------------------------|-----------------------|--------|---------------|--|
| Name | Date Modified | Size | Kind | |
| antlr-2.7.6.jar | 4 January 2011, 15:53 | 404 KB | Java JAR file | |
| commons-pool-1.2.jar | 4 January 2011, 15:53 | 48 KB | Java JAR file | |
| COPYING-libs | 4 January 2011, 15:53 | 68 KB | Plain text | |
| exist-modules.jar | 4 January 2011, 15:53 | 72 KB | Java JAR file | |
| exist.jar | 4 January 2011, 15:53 | 3.3 MB | Java JAR file | |
| idx.jar | 4 January 2011, 15:53 | 144 KB | Java JAR file | |
| je.jar | 4 January 2011, 15:53 | 2 MB | Java JAR file | |
| jgroups-all.jar | 4 January 2011, 15:53 | 1.6 MB | Java JAR file | |
| jung.jar | 4 January 2011, 15:53 | 956 KB | Java JAR file | |
| log4j-1.2.14.jar | 4 January 2011, 15:53 | 392 KB | Java JAR file | |
| MinML2.jar | 4 January 2011, 15:53 | 20 KB | Java JAR file | |
| prefuse.jar | 4 January 2011, 15:53 | 564 KB | Java JAR file | |
| README | 4 January 2011, 15:53 | 8 KB | Plain text | |
| resolver.jar | 4 January 2011, 15:53 | 68 KB | Java JAR file | |
| sunxacml.jar | 4 January 2011, 15:53 | 264 KB | Java JAR file | |
| teccli.jar | 4 January 2011, 15:53 | 216 KB | Java JAR file | |
| teccli.properties | 4 January 2011, 15:53 | 4 KB | Document | |
| tecclipluginlist.txt | 4 January 2011, 15:53 | 4 KB | Plain text | |
| xmlldb.jar | 4 January 2011, 15:53 | 16 KB | Java JAR file | |
| xmlrpc-1.2-patched.jar | 4 January 2011, 15:53 | 116 KB | Java JAR file | |

- Click on **teccli.jar**(highlighted above). (NB. It might appear as just **teccli**.)
- A window will appear asking you to select a corpus. The default corpus is TEC, which is available on-line. To view your own corpus, select “Choose new local corpus”.
- A window will open that asks you to “Choose location for the index”. Select the folder that the index is saved in, i.e. the **epi** (or **jpi**) folder.
- A window may open that asks you to “Choose a headers directory”. Select the folder that the header files are saved in, i.e. the **eph** (or **jph**) folder.
- The concordance browser window will appear.
- Type “*this*” to check the concordancer is working.

Troubleshooting

If you have followed these instructions but the concordance browser does not launch, try the following:

- Check what versions of Java you have installed.

For Mac OS X users:

Open Utilities > Java Preferences> General tab.

Move the most recent version to the top of the list.

- Check which version of Java you are using.

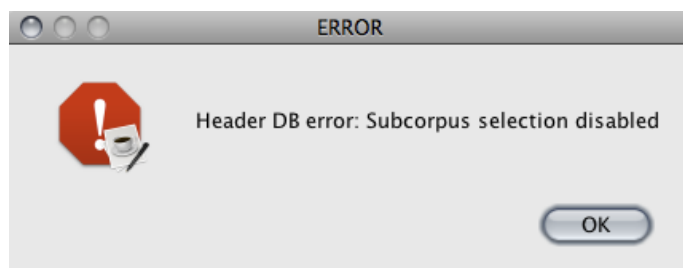
For Mac OS X users:

Open Terminal(this can be found in Applications > Utilities > Terminal)

Type “`java -version`”, then Enter

- Confirm that this version is compatible with the current version of the Corpus Browser.
- Check your Java security settings to confirm that you are not blocking anything (Utilities > Java Preferences > Security tab).

If the concordance browser launches but gives this error message:



Click **OK** and proceed.

3.6 Sharing your corpus

If you want to make your corpus available to other users, you will need to download the Corpus Server module (modnlp-tecser).

4. Searching a Corpus

4.1 The Corpus Browser Interface

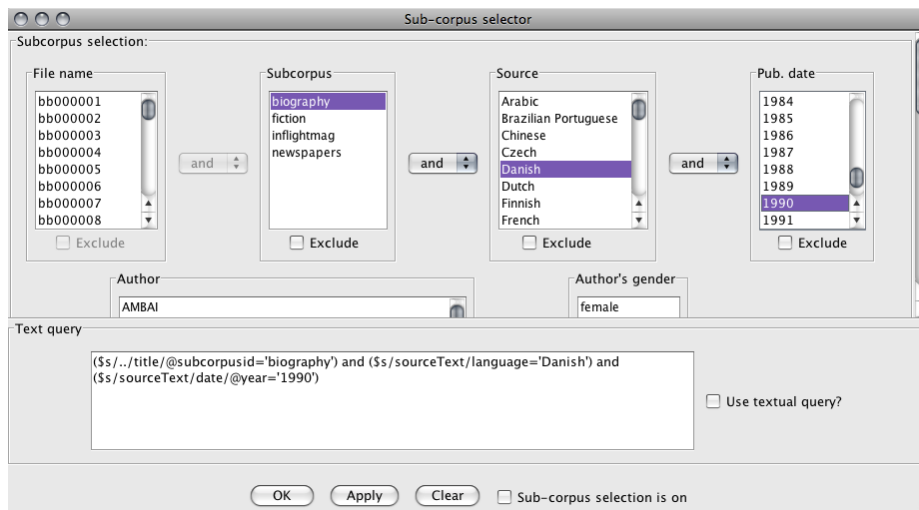
Many of the features of the Corpus Browser are self-explanatory but brief descriptions of some of the menu options are provided below.

File Menu

- **New local corpus**
If you have saved a corpus on your computer, you can select it by selecting the folder where the index is stored, e.g. the **epi** folder.
- **New Internet corpus**
If you want to access a corpus remotely you should input the IP address of the server on which the corpus is stored. The default IP address provides access to the TEC internet corpus.
- **Save Concordances**
You can save the results of your corpus searches. To do this, click “Save Concordances” > specify a filename and location > “Save to Disc”. You can open the saved file with a text editor or Word. To set out the concordances clearly in Word, go to Page Setup > Orientation > Landscape. Then set the font to Courier size 7.5.

Options Menu

- **Case sensitive**
The default setting for corpus searches is that they are not case-sensitive. If you want to carry out case-sensitive searches, check this option.
- **Select Sub-corpus**
This option opens the “Select Sub-corpus” window. In this window you can select particular text files or groups of text files to search. There are a number of selection options including date of publication, gender of translator, source text language etc. The options available are based on the information stored in the header file for each text in the corpus.

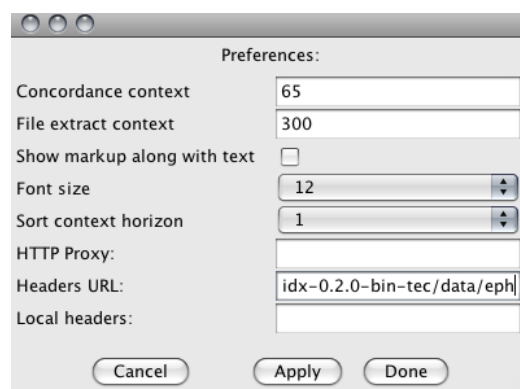


- **Activate sub-corpus selection**

This option activates and deactivates the selection criteria set in the “Select Sub-corpus” menu. When this option is deselected, any searches will be carried out on the whole corpus.

Preferences

- This option opens the “Preferences” window in which you can change various settings.



Some settings you may wish to experiment with include:

Concordance context> changes the number of words that are shown in concordance lines.

File extract context> changes the number of words that are shown when using the “Extract” function to check individual concordances in their broader context.

Show markup along with text> shows you any xml markup that is present in the text file. For example, </p> indicating a paragraph break.

Plugins Menu

- **Word frequency list**

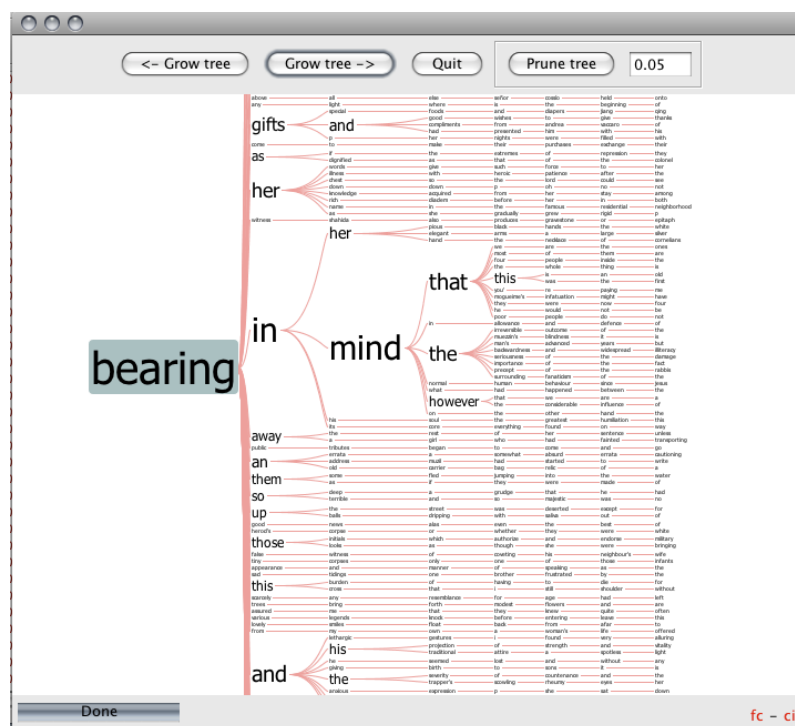
This tool generates a frequency list for the corpus. You can change the number of words listed (e.g. 100 most frequent words, or 500 most frequent words) and click “Get List”. You can save lists which can later be opened using a text editor or Word.

- **Corpus description browser**

This lists the filenames, title and other information (meta-data) for the texts in the corpus. The number of tokens and type/token ratio for each text are also listed.

- **Concordance tree viewer**

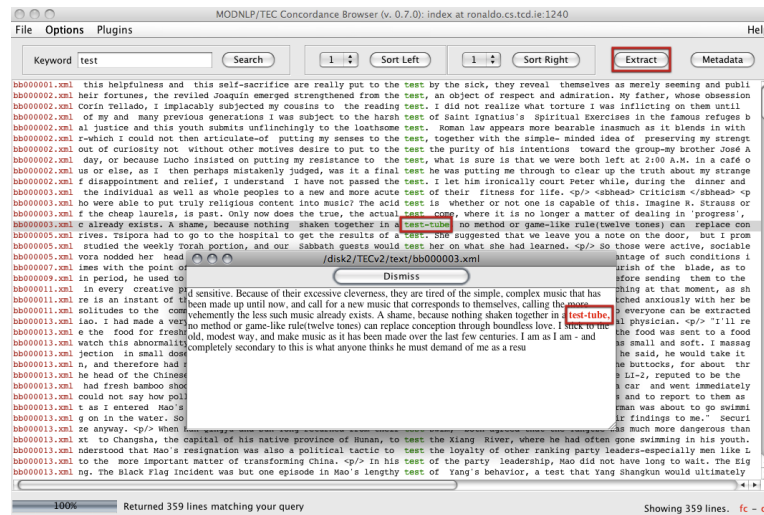
This is a corpus visualization tool that enables you to “grow” a tree for various keywords. The tree view can help you to identify frequently occurring collocations, since the size of the text reflects frequency of occurrence in the corpus.



Other Tools

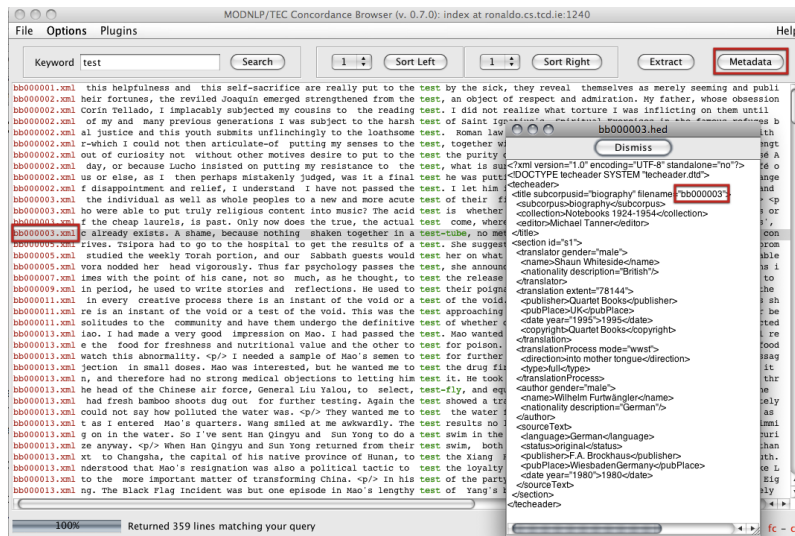
- **Extract**

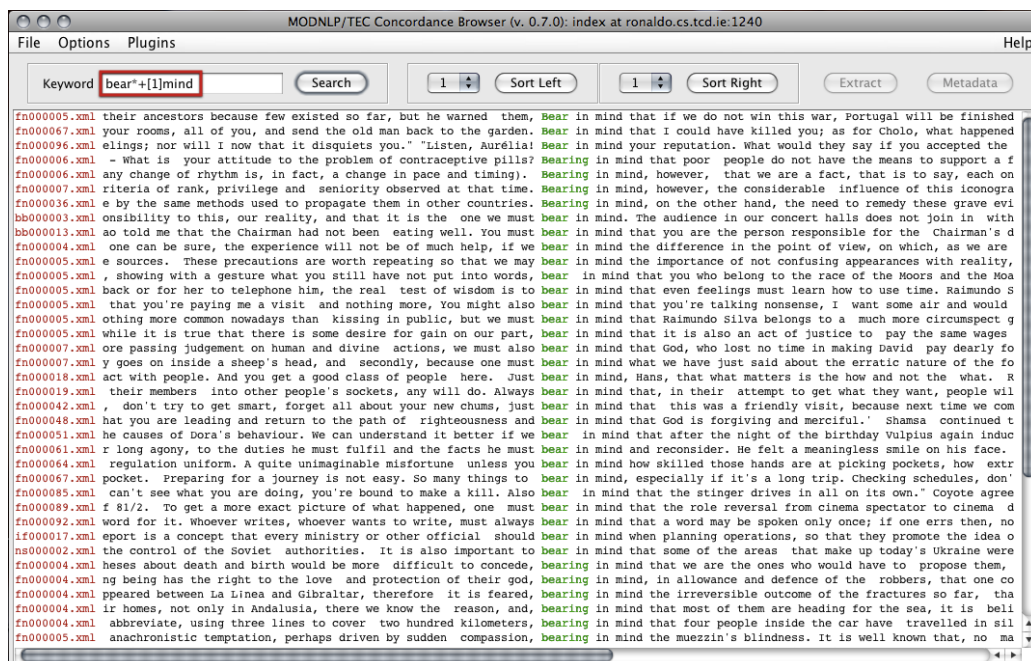
Allows you to see more context for a particular concordance. Select the concordance line by clicking on it and then “Extract”.



- **Metadata**

Allows you to see the information stored in the header file for a particular concordance. Select the concordance line by clicking on it and then “Metadata”.





4.3 Sorting Concordances

When you have generated concordance lines, you can arrange them in various ways to help you identify patterns. There are three ways to do this.

- **Sort Tool**
Allows you to sort concordances to the left or right, and specify words.
- **Tree Viewer**
- **Mosaic Viewer**

4.4 Saving Search Results

You can save concordances, frequency lists, etc for use later. In any window click on **Save**. The saved file contains text which can be opened using a Text Editor or Word. To set out concordances clearly in Word:

- Go to Page Setup > Orientation > Landscape.
- Set the font to Courier size 7.5.

5. Other Resources

Other resources you may wish to consult include:

<http://www.monabaker.com/tsresources/TranslationalEnglishCorpus.htm>

This website provides information about TEC and its contents, a select bibliography, and a PDF presentation by Dr. Saturnino Luz, introducing the TEC Tools corpus software.

<http://ronaldo.cs.tcd.ie/tec2/jnlp/>

This website provides access to the TEC resource using the Corpus Browser. You can use this to familiarize yourself with the Browser interface and try out searches as described in section 4.4.

There are a number of publications that describe various aspects of corpus linguistics that you may find useful. In particular, the following volume, due for publication in September 2011, talks specifically about the TEC Tools (MODNLP) software and the XML encoding it uses.

6. References

Luz, S. (2011) 'Web-based Corpus Software' in A. Kruger, K. Wallmach and J. Munday (eds) *Corpus-Based Translation Studies: Research and Applications*, chapter 5, pages 124-149. London: Continuum.